

uninformed search ①

Frontier = {A, D, B, E, H, C, F, G, I, I₄₆}

explored = {A, D, B, C, E, H, G, I}

total cost = 46

Frontier = {A}, explored = {}

توضیح مرحله به مرحله : Frontier ← A ①

Frontier = {B, D}, explored = {A}

Frontier ← D, B / explored ← A ②

Frontier = {B, E}, explored = {A, D}

Frontier ← E / explored ← D ③
A → D

Frontier = {E, H, C}, explored = {A, D, B}

Frontier ← H, C / explored ← B ④
A → B

Frontier = {E, H, F, G}, explored = {A, D, B, C}

Frontier ← G, F / explored ← C ⑤
A → B → C

Frontier = {H, F, G}, explored = {A, D, B, C, E}

Frontier ← F / explored ← E ⑥
delete F₅₀ , A → D → E

Frontier = {F, I₄₆, G}, explored = {A, D, B, C, E, H}

Frontier ← I₄₆ / explored ← H ⑦
A → B → H

Frontier = {F, I₄₆}, explored = {A, D, B, C, E, H, G}

Frontier ← X / explored ← G ⑧
A → B → C → G

cost 46

✓ / explored ← I ⑨
A → B → H → I

A → B → H → I

هیورستیک پیشنهادی این است که فاصله‌ی هر عنصر با جایگاه درست خودش را از لحاظ میزان قدر مطلق ارتفاع بدست آورده روی تمامی node ها جمع بپذیم و تقسیم بر ۲ کنیم.

مثبت تقسیم بر ۲:

در صورتی که شکل مقابل در دسترس باشد:

در دو مطابق (a) باشد با یکبار جایابی a, b در شرایط واقعی (cost) تا هزینه

در دو (ما هیورستیک ما قبل از تقسیم بر ۲ مقدار ۲ نشان داده است پس باید تقسیم بر ۲ کنیم (بدونین حالت این است)

در واقع در جایابی هر ۲ عنصر و نحوه ۲ حالت وجود دارد:

۱) جایابی a, b (cost = 1) باعث شود a یکی در شود و b یکی نزدیک شود (یعنی یک ارتفاع)

یعنی h ما - پس نزدیک شدن a یکی کم شود - پس دور شدن b یکی اضافه شده است
یعنی h تغییر نمی کند (هیورستیک ثابت است) = مقدار تغییر h بین دو state از cost کمتر است

۲) جایابی a, b (cost = 2) باعث شود هر ۲ در جایگاه اصلی در شوند یا هر ۲ به جایگاه اصلی

نزدیک شوند در این حالت مجموع فاصله‌ی ارتفاع ها به اندازه ۲ تا زیاد می شود که پس از تقسیم بر ۲

- همان عدد ۱ که مساوی cost است میسیم. در نتیجه تفاوت cost های بین ۲ state

ب جایابی (یعنی ۲ عنصر ۱ است و تفاوت h های آن ها نزدیک است. «تفاوت نزدیک تر»

تفاوت h های دو state صفر و cost بین آنها ۱ است که «این حالت نیز تفاوت بین state h

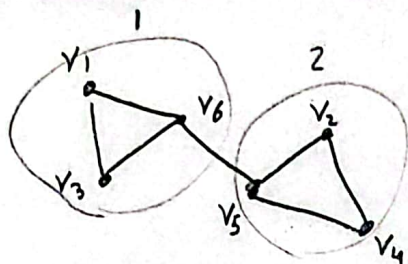
۲ state از cost بین ۲ state کمتر است. یعنی consistent با heuristic است.

پس در صورتی که یک هیورستیک consistent باشد \Rightarrow admissible نیز هست پس هیورستیک ما

هم admissible و هم consistent است.

$$\rightarrow h = \frac{\sum |height(i_t) - height(i_{t+1})|}{2}$$

الف) هر کدوم شامل n زن است که هر زن نمایندگی یک node است. شامل شماره‌های گره‌ها.
 node است که به صورت Random خنثی می‌شود. شماره‌های گره‌ها می‌تواند هر عدد رندمی بین n تا n (نماد node ها است)



1	2	1	2	2	1
v_1	v_2	v_3	v_4	v_5	v_6

ب)

ج) هر زن با احتمال P_m با یک شماره گره جدید در باره n جایگزینی کنیم.

د) هر زن از هر کدوم را با احتمال P_c با همان عنصر (زن) از کدوم دیگر جایگزینی کنیم.

ه) یکی از Fitness Function های مناسب برای هر کدوم استفاده از شاخص silhouette است. این شاخص برای یک دسته‌بندی (clustering) معین نسبت خوشه‌بندی را می‌سنجد.

$$a(i) = \frac{1}{|C_i| - 1} \sum_{j \in C_i, j \neq i} d(i, j)$$

که C_i یعنی خوشه‌ای که node شماره i را به آن اختصاص داده‌ایم و $d(i, j)$ در واقع فاصله‌ی node های i و j است. در صورتی که بین i و j یک مقدار یک و در غیر این صورت مقدار صفر می‌گیرد. در نهایت ما می‌توانیم $a(i)$ را بعنوان (شبه حد متوسط) خوشه خود اختصاص دارد در نظر بگیریم. حال باید متوجه شویم که نقطه i باید به خوشه‌ای باشد C_k - عنوان می‌کنیم فاصله از تمامی نقاط داخل C_k تعیین می‌کنیم.

$$b(i) = \min_{k \neq i} \frac{1}{|C_k|} \sum_{j \in C_k} d(i, j)$$

آن خوشه که در \min خارج شده است یعنی خوشه‌ای که i را به آن اختصاص داده‌ایم (بهترین شاخص یا محل قرارگیری) در میان خوشه‌ها دارد. حال مقدار سیلوئت را از روش زیر به دست می‌آوریم:

$$S(i) = \frac{b(i) - a(i)}{\max(a(i), b(i))}$$

در صورتی که $a = 1$ یا $c = 1$ باشد $\Rightarrow a$ یا c است در این حالت

$$s(i) = \begin{cases} 1 - a(i)/b(i) & \text{if } a(i) < b(i) \\ 0 & \text{if } a(i) = b(i) \\ b(i)/a(i) - 1 & \text{if } a(i) > b(i) \end{cases}$$

$$1 \leq s(i) \leq 1$$

از ترفیع با درشتن می شود

حالت: ترفیع، Fitness Function می نامیم، به برای حرکت در تمام مجموعه $s(i)$ های زن های آن را Fitness Func تعریف می کنیم.

چون $s(i)$ ترفیع یعنی $a(i) < b(i)$ است و از آنجایی که $a(i)$ یک معیار برای چگونگی عدم شباهت نا با خود خورش است و بزرگ بودن $a(i)$ یعنی نا با خودهای مشابه تطابق خوبی نداشته است. پس $s(i)$ ترفیع $\Rightarrow 1$ یعنی خود شباهتی مناسب نا است.

در مقابل آن $s(i)$ به معنی یک ترفیع با هم با همان منطق بالا منجر می شویم که آن نا در خود شباهتی بود (خود شباهت) تطابق بیشتری است.

پس زیاد شدن مجموع $s(i)$ ها یعنی یک شباهتی خوب و کم بودن آن یعنی یک شباهتی بد.

که این می تواند یک شاخص Fitness مناسب برای مسئله باشد.

	s_0	s_1	s_2	s_3
π_1	a_1	a_0	a_0	a_0
π_2	a_2	a_0	a_0	a_0

در این بازی ۲ پالیسی داریم.

$$V^*(s_0) = \max_a 0 + \gamma \sum_{s'} T(s, a, s') V^*(s') = \gamma \max \{V^*(s_1), V^*(s_2)\}$$

$$V^*(s_1) = \max_a 0 + \gamma \sum_{s'} T(s, a, s') V^*(s') = \gamma [(1-p)V^*(s_1) + pV^*(s_3)]$$

$$V^*(s_2) = \max_a 1 + \gamma \sum_{s'} T(s, a, s') V^*(s') = 1 + \gamma [(1-q)V^*(s_0) + qV^*(s_3)]$$

$$V^*(s_3) = \max_a 10 + \gamma \sum_{s'} T(s, a, s') V^*(s') = 10 + \gamma V^*(s_0)$$

توجه کنید (یعنی فرض کنیم $p=0$)

$$\pi^*(s_0) = \arg \max_a \gamma \sum_{s'} T(s, a, s') V^*(s') \rightarrow \text{بسیار } \gamma \neq 0$$

در صورت بودن γ هیچ π یکسانی به دست نمی آید.

پس در صورت $p=0$ ، $V^*(s_2) > V^*(s_1)$ باشد، $\pi^*(s_0) = a_2$

حال p را بررسی می کنیم. در صورتی که $p=0$ باشد چون $V^*(s_1)$ را بررسی می کرد پس در صورت $p=0$ (۱۶۹)

$$V^*(s_2) = 1 + \gamma [(1-q)V^*(s_0) + qV^*(s_3)] \geq 1$$

$$V^*(s_2) > V^*(s_1)$$

پس

و می بینیم که چون $T(s, a_2, s')$ برابر است پس

$$\pi^*(s_0) = a_2$$

می شود چرا که طبق فرمول $\pi^*(s_0)$ بین $V^*(s_2)$ و $V^*(s_1)$ ما $V^*(s_2)$ را انتخاب می کنیم.

پس برای $p=0$ فرض می کردیم که $\pi^*(s_0) = a_2$ است چرا که برای a_1 فقط هیچ سودی ندارد.

این مسئله دقیقاً برعکس مسئله قبل است.
 جواب منفی است؟ خیر نمی‌دانیم.

$$V^*(s_2) = 1 + \gamma [(1-p)V^*(s_0) + pV^*(s_3)] \geq 1$$

چرا که

نیست اگر نخواهیم در هر صورت در state s_0 action a_0 را انتخاب کنیم پس باید از آنجا

$$V^*(s_0) > V^*(s_2)$$

بایست. می‌دانیم $V^*(s_2) \geq 1$ است. چون

$$V^*(s_0) = \gamma [(1-p)V^*(s_0) + pV^*(s_3)] = \frac{\gamma p V^*(s_3)}{1 - \gamma(1-p)}$$

حال حول مباحث محدودیتی برای γ نداریم، مگر می‌توانیم یک γ پیدا کنیم که $V^*(s_0) < 1$ باشد و به این دلیل شرط رسم برای آنکه $V^*(s_0) > V^*(s_2)$ باشد برقرار نیست (یعنی نقض یک γ است) که این شرط را نقض کند پس هیچ محدودیتی برای γ نداریم که با آن از $\gamma = 0.9$ بزرگتر باشد.

$$\gamma = 0.9 \quad p = 0.25 \quad \epsilon = 10^{-3} \checkmark$$

$$V^*(s_0) = 14.184$$

$$\pi^*(s_0) = a_1$$

جواب نهایی این مسئله

$$V^*(s_1) = 15.760$$

$$\pi^*(s_1) = a_0$$

$$V^*(s_2) = 15.696$$

$$\pi^*(s_2) = a_0$$

$$V^*(s_3) = 22.766$$

$$\pi^*(s_3) = a_0$$

$$\begin{aligned} V^*(s_0) &= \max_a \left\{ \left[1 \times (0 + \gamma V^*(s_1)) \right], \left[1 \times (0 + \gamma V^*(s_2)) \right] \right\} = 0.9 \\ V^*(s_1) &= \max_a \left\{ \left[0.25 \times (0 + \gamma V^*(s_0)) + 0.75 \times (1 + \gamma V^*(s_1)) \right], \left[0.25 \times (0 + \gamma V^*(s_0)) + 0.75 \times (1 + \gamma V^*(s_2)) \right] \right\} = 3.8575 \\ V^*(s_2) &= \max_a \left\{ \left[0.25 \times (0 + \gamma V^*(s_0)) + 0.75 \times (1 + \gamma V^*(s_1)) \right], \left[0.25 \times (0 + \gamma V^*(s_0)) + 0.75 \times (1 + \gamma V^*(s_2)) \right] \right\} = 2.25 \\ V^*(s_3) &= \max_a \left\{ 1 \times (10 + \gamma V^*(s_3)) \right\} = 10.81 \end{aligned}$$

iteration 1

همین ترتیب را ادامه می‌دهیم.

این action قبل از این است که به مقادیر عددی درستی دست یابیم و آنقدر محاسبه می‌کنیم

$$V^*(s) = \max_a \sum_{s'} T(s, a, s') (R + \gamma V^*(s'))$$

$$V^*(s_0) = \max_a \{ [1 \times (0 + 0.25 \times 2.5)], [1 \times (0 + 0.9 \times 3.8575)] \} = 3.47175$$

$$V^*(s_1) = \max_a T(s, a, s') [R(s, a, s') + \gamma V^*(s')] = 3.951$$

$$V^*(s_2) = \text{---} = 7.461$$

$$V^*(s_3) = \text{---} = 14.678$$

به همین ترتیب ادامه می دهیم تا زمانی که میزان خطای هر ایتِ رِشین یا در واقع تفاوت ویدیو ویدیو ایتِ رِشین قبلی کمتر از ϵ که همان 0.1 است شود بین کانوِج شده است و مفاد در V به دست آمده مفاد در درست هستند. همچنین آن a که باعث V^* می شود نیز π^* نامیده می شود.

که در این سؤال:

$$V^*(s_0) = 14.184 \quad \pi^*(s_0) = a_1$$

$$V^*(s_1) = 15.76 \quad \pi^*(s_1) = a_0$$

$$V^*(s_2) = 15.696 \quad \pi^*(s_2) = a_0$$

$$V^*(s_3) = 22.766 \quad \pi^*(s_3) = a_0$$

بدان 53 تا سوال

سه توضیح درست آوردن 1, 1, 1

$$V^*(s_0) = 0, \quad V^*(s_1) = 0$$

$$V^*(s_2) = 0.25 \times (1 + 0.9 \times 0) + 0.75 \times (1 + 0.9 \times 0) = 1$$

$$V^*(s_3) = 0.25 \times (10 + 0.9 \times 0) + 0.75 \times (10 + 0.9 \times 0) = 10$$

سوال دوم

۱۱ قضای حالت: چون در حضور مسافران مطلع نیستیم پس قضای حالت در ناگسی و ندری
آن است.

A	B
C	D

ناگسی در A است
4
X
2
ناگسی مسافر در B است

۸ حالت

آنکه جایگاه مسافران بیرونی باشد و این با احتمال برای وجود مسافر در این
آن پس در جایگاه مسافر مسافر کردن مسافر را اعلام می کنیم. در صورتی که مسافر مسافر بود

۱۲ ما ۷ تا action داریم. ۱۱ مسافر ماندن ۱۲ مسافر کردن ۱۳ مسافر کردن ۱۴ مسافر کردن ۱۵ مسافر کردن ۱۶ مسافر کردن ۱۷ مسافر کردن

۱۳ در صورتی که ناگسی مسافر در آن باشد و در مسافر کردن مسافر در آن باشد به خود آن است و بر آن مسافر (تغییر کند)

در حالتی که ناگسی مسافر در آن باشد و در مسافر کردن مسافر در آن باشد به خود آن است و بر آن مسافر (تغییر کند)
مسافر مسافر در آن باشد و در مسافر کردن مسافر در آن باشد به خود آن است و بر آن مسافر (تغییر کند)

در حالتی که ناگسی مسافر در آن باشد و در مسافر کردن مسافر در آن باشد به خود آن است و بر آن مسافر (تغییر کند)
در حالتی که ناگسی مسافر در آن باشد و در مسافر کردن مسافر در آن باشد به خود آن است و بر آن مسافر (تغییر کند)
جایگاه در مسافر مسافر می بینیم

در صورت انجام حرکت از مسافر است و به احتمال مسافر در آن باشد به خود آن است و بر آن مسافر (تغییر کند)

در مسافر مسافر ماندن هم در حال State می بینیم.

در تمامی حالات به احتمال مسافر است با انجام یک مسافر آن مسافر در آن باشد و در مسافر کردن مسافر در آن باشد
به خود آن است و بر آن مسافر (تغییر کند)

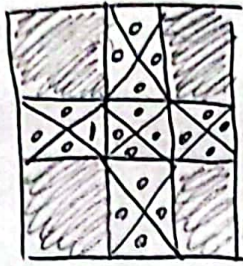
۴- ما برای هر عمل مانند سود کردن، پیاوردن و ... یک Reward دریافت می‌کنیم.

در واقع هر بار جویندگی ما شدن در محیط (انتظار) و جویندگی (انتظار) پیاوردن هر خانه، جویندگی سود کردن مسافر هر خانه درست، جویندگی پیاوردن در خانه درست و جویندگی پیاوردن در خانه اشتباه را باید معینان value برای state در نظر گرفت. حال متوجه شدیم که ما یک جویندگی را داشتیم که محیط صفر می‌شود یا در صورتی که سازی نداشته باشیم جویندگی پیاوردن صفر می‌شود. با وجود حالات مختلف مجموع این جویندگی و جویندگی سود کردن و معینان Reward - agent را در می‌آوریم.

در ضمن برای رسیدن به هدف صورت سوال باید جویندگی را داشتیم که جویندگی شدت ماندن بیشتر و جویندگی پیاوردن اشتباه بیشتر و همین جویندگی پیاوردن درست از مابقی بیشتر باشد.

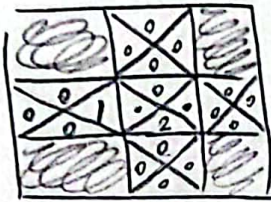
در نهایت با محاسبی تابع جویندگی و جویندگی‌ها برای هر action q-value تعیین کردن اما Converge بیش می‌آوریم و بیشترین q-value را انتخاب می‌کنیم معینان π^* در نظر می‌آوریم.

$$Q(s,a) = (1-\alpha)Q(s,a) + \alpha [R(s,a,s') + \gamma \max_{a'} Q(s',a')] \quad \text{RL 0}$$



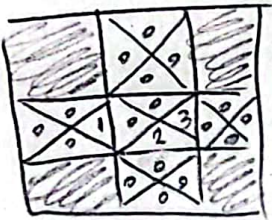
$$(B, East, C, 2) - 1$$

$$(1-0.5) \times 0 + 0.5 [2 + 1 \times 0] = 1$$



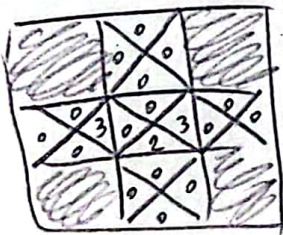
$$(C, South, E, 4) - 2$$

$$(1-0.5) \times 0 + 0.5 (4 + 1 \times 0) = 2$$



$$(C, East, A, 6) - 3$$

$$(1-0.5) \times 0 + 0.5 \times (6 + 1 \times 0) = 3$$



$$(B, East, C, 2) - 3$$

$$(1-0.5) \times 1 + 0.5 \times (2 + 1 \times 3) = 3$$

→ q-table ✓

(ب) در حالتی که ما از روش مذکور استفاده کنیم در واقع (یعنی) درست را داریم و تنها باید ۷ بار پیدا کنیم.

در اینجا با احتمال ۴-۱ روی (یعنی) درست $it \leq 10$ می‌کنیم و در صدی که با احتمال ۴ نیز یک

action نادرست داشته باشیم در تعداد iteration های بالا باز ۷۵٪ های درست را استخراج می‌کنیم. و چون

در حال استفاده از ϵ -greedy هستیم پس در صورت کاهش ϵ نیز در اصل بالاتر از ۷۵٪ از ۷۵٪ های درست را می‌گیریم.

مقایسه (۱) Performance: به دلیل داشتن جریب درست این روش سریعتر اجرا می‌شود.

۱۲ در اکثریت Q -learning معمولی هرگز خطا در تقریب هنگام به روز رسانی Q -table قابل

اصلاح است اما در این روش این امکان بالقوه کمتر است اما میزان خطای بسیار کمتر است.

(۱۳) معیار مقایسه: می‌توان گفت که به دلیل اینکه هر بار در حال به روز کردن همان Q -value

که آن در (یعنی) پیدا می‌کنیم حتمی این مقدار سریعتر و دقیق‌تر به روز رسانی می‌شوند.

به طور کلی معیار نتیجه‌گیری می‌توان گفت وقتی از (یعنی) بهینه‌ی پیرامونی می‌گیریم عملکرد اکثریت بهینه‌ی پیرامونی

کمتر از عملکرد Q -learning دارد است. زیرا احتمال انتخاب عمل اشتباه در آن ۴ است.